

Optimizing Ranking Measures for Compact Binary Code Learning^{*}

Guosheng Lin¹, Chunhua Shen^{1**}, Jianxin Wu²

¹University of Adelaide, Australia

²Nanjing University, China

Abstract. Hashing has proven a valuable tool for large-scale information retrieval. Despite much success, existing hashing methods optimize over simple objectives such as the reconstruction error or graph Laplacian related loss functions, instead of the performance evaluation criteria of interest—multivariate performance measures such as the AUC and NDCG. Here we present a general framework (termed StructHash) that allows one to directly optimize multivariate performance measures. The resulting optimization problem can involve exponentially or infinitely many variables and constraints, which is more challenging than standard structured output learning. To solve the StructHash optimization problem, we use a combination of column generation and cutting-plane techniques. We demonstrate the generality of StructHash by applying it to ranking prediction and image retrieval, and show that it outperforms a few state-of-the-art hashing methods.

1 Introduction

The ever increasing volumes of imagery available, and the benefits reaped through the interrogation of large image datasets, have increased enthusiasm for large-scale approaches to vision. One of the simplest, and most effective means of improving the scale and efficiency of an application has been to use hashing to pre-process the data [10,23,14,19,16,13].

Depending on applications, specific measures are used to evaluate the performance of the generated hash codes. For example, information retrieval and ranking criteria [17] such as the Area Under the ROC Curve (AUC) [7], Normalized Discounted Cumulative Gain (NDCG) [6], Precision-at-K, Precision-Recall and Mean Average Precision (mAP) have been widely adopted to evaluate the success of hashing methods. However, to date, most hashing methods are usually learned by optimizing simple errors such as the reconstruction error (e.g., binary reconstruction embedding hashing [10]) or the graph Laplacian related loss [25,16,23]. These hashing methods construct a set of hash functions that map the original high-dimensional data into a much smaller binary space, typically with the goal of preserving neighborhood relations. The resulting compact binary encoding enables fast similarity computation between data points by using the

^{*} Appearing in Proc. European Conf. Computer Vision, 2014.

^{**} Corresponding author (e-mail: chunhua.shen@adelaide.edu.au).

Hamming distance, which can be carried out by rapid, often hardware-supported, bit-wise operations. Furthermore, compact binary codes are much more efficient for large-scale data storage.

To our knowledge, none of the existing hashing methods have tried to learn hash codes that *directly* optimizes a multivariate performance criterion. In this work, we seek to reduce the discrepancy between existing learning criteria and the evaluation criteria (such as retrieval quality measures).

The proposed framework accommodates various complex multivariate measures. By observing that the hash codes learning problem is essentially an information retrieval problem, various ranking loss functions can and should be applied, rather than merely pairwise distance comparisons. This framework also allows to introduce more general definitions of “similarity” to hashing beyond existing ones.

In summary, our main contributions are as follows.

1. We propose a flexible binary hash codes learning framework that directly optimizes complex multivariate measures. This framework, *for the first time*, exploits the gains made in structured output learning for the purposes of hashing. Our hashing method, labelled as StructHash, is able to directly optimize various multivariate evaluation criteria, such as information retrieval measures (e.g., AUC and NDCG [6]).
2. To facilitate StructHash, we combine column generation and cutting-plane methods in order to efficiently solve the resulting optimization problem, which may involve exponentially or even infinitely many variables and constraints.
3. Applied to ranking prediction for image retrieval, the proposed method demonstrates state-of-the-art performance on hash function learning.

2 Related work

One of the best known *data-independent* hashing methods is locality sensitive hashing (LSH) [3], which uses random projection to generate binary codes. Recently, a number of *data-dependent* hashing methods have been proposed. For example, spectral hashing (SPH) [23] aims to preserve the neighbourhood relation by optimizing the Laplacian affinity. Anchor graph hashing (AGH) [16] makes the original SPH much more scalable. Examples of supervised or semi-supervised hashing methods include binary reconstruction embedding (BRE) [10], which aims to minimise the expected distances; and the semi-supervised sequential projection learning hashing (SPLH) [21], which enforces the smoothness of similar data points and the separability of dissimilar data points.

To obtain a richer representation, kernelized LSH [11] was proposed, which randomly samples training data as support vectors, and randomly draws the dual coefficients from a Gaussian distribution. Liu et al. extended Kulis and Grauman’s work to kernelized supervised hashing (KSH) [15] by learning the dual coefficients instead. Lin et al. [13] employed ensembles of decision trees as the hash functions. Nonetheless, all of these methods do not directly optimize the

multivariate performance measures of interest. We formulate hash codes learning as a structured output learning problem, in order to directly optimize a wide variety of evaluation measures.

We are primarily inspired by recent advances in learning to rank such as the metric learning method in [17], which directly optimizes several different ranking measures. We aim to learn hash functions, which leads to a very different learning task preventing directly applying techniques in [17]. We are also inspired by the recent column generation based hashing method, column generation hashing (CGH) [12], which iteratively learns hash functions using column generation. However, their method optimizes the conventional classification-related loss, which is much simpler than the multivariate loss that we are interested in here. Moreover, the optimization of CGH relies on all triplet constraints while our method is able to use much less number of constraints without sacrificing the performance.

Our framework is built on the structured SVM [20], which has been applied to many applications for complex structured output prediction, e.g., image segmentation, action recognition and so on.

Notation Let $\{(\mathbf{x}_i; \mathbf{y}_i)\}$, $i = 1, 2, \dots$, denote a set of input and output pairs. The discriminative function for structured output prediction is¹ $F(\mathbf{x}, \mathbf{y}) : \mathcal{X} \times \mathcal{Y} \mapsto \mathbb{R}$, which measures the compatibility of the input and output pair (\mathbf{x}, \mathbf{y}) . Given a query \mathbf{x}_i , we use \mathcal{X}_i^+ and \mathcal{X}_i^- to denote the subsets of relevant and irrelevant data points in the training data. Given two data points: \mathbf{x}_i and \mathbf{x}_j , $\mathbf{x}_i \prec_{\mathbf{y}} \mathbf{x}_j$ ($\mathbf{x}_i \succ_{\mathbf{y}} \mathbf{x}_j$) means that \mathbf{x}_i is placed before (after) \mathbf{x}_j in the ranking \mathbf{y} .

2.1 Structured SVM

First we provide a brief overview on structured SVM. Structured SVM enforces that the score of the “correct” model \mathbf{y} should be larger than all other “incorrect” model \mathbf{y}' , $\forall \mathbf{y}' \neq \mathbf{y}$, which writes:

$$\forall \mathbf{y}' \in \mathcal{Y} : \quad \mathbf{w}^\top [\psi(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{x}, \mathbf{y}')] \geq \Delta(\mathbf{y}, \mathbf{y}') - \xi. \quad (1)$$

Here ξ is a slack variable (soft margin) corresponding to the hinge loss. $\psi(\mathbf{x}, \mathbf{y})$ is a vector-valued joint feature mapping. It plays a key role in structured learning and specifies the relationship between an input \mathbf{x} and output \mathbf{y} . \mathbf{w} is the model parameter. The label loss $\Delta(\mathbf{y}, \mathbf{y}') \in \mathbb{R}$ measures the discrepancy of the predicted \mathbf{y}' and the true label \mathbf{y} . A typical assumption is that $\Delta(\mathbf{y}, \mathbf{y}) = 0$, $\Delta(\mathbf{y}, \mathbf{y}') > 0$ for any $\mathbf{y}' \neq \mathbf{y}$, and $\Delta(\mathbf{y}, \mathbf{y}')$ is upper bounded. The prediction \mathbf{y}^* of an input \mathbf{x} is achieved by

$$\mathbf{y}^* = \operatorname{argmax}_{\mathbf{y}} F(\mathbf{x}, \mathbf{y}) = \mathbf{w}^\top \psi(\mathbf{x}, \mathbf{y}). \quad (2)$$

¹ To be precise, the prediction function should be written as $F(\mathbf{x}, \mathbf{y}; \mathbf{w})$ because it is parameterized by \mathbf{w} . For simplicity, we omit \mathbf{w} .

Optimization for the model parameter \mathbf{w} For structured problems, the size of the output $|\mathcal{Y}|$ is typically very large. Considering all possible constraints in (1) is generally intractable. The cutting-plane method [9] is commonly employed, which allows to maintain a small working-set of constraints and obtain an approximate solution of the original problem up to a pre-set precision. To speed up, a 1-slack reformulation is proposed [8]. Nonetheless the cutting-plane method needs to find the most violated label (equivalent to an inference problem)

$$\operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}} \mathbf{w}^\top \psi(\mathbf{x}, \mathbf{y}') + \Delta(\mathbf{y}, \mathbf{y}'). \quad (3)$$

Structured SVM typically requires: 1) a well-designed feature representation $\psi(\cdot, \cdot)$; 2) an appropriate label loss $\Delta(\cdot, \cdot)$; 3) solving inference problems (2) and (3) efficiently.

Ranking prediction with structured output In a retrieval system, given a test data point \mathbf{x} , the goal is to predict a ranking of data points in the database. For a “correct” ranking, relevant data points are expected to be placed in front of irrelevant data points. A ranking output is denoted by \mathbf{y} . Let us introduce a symbol $y_{jk} = 1$ if $\mathbf{x}_j \prec_{\mathbf{y}} \mathbf{x}_k$ and $y_{jk} = -1$ if $\mathbf{x}_j \succ_{\mathbf{y}} \mathbf{x}_k$. The ranking can be evaluated by various measures e.g. AUC, NDCG, mAP. These evaluation measures can be optimized directly as label losses Δ [7,17]. Here $\psi(\mathbf{x}, \mathbf{y})$ can be defined as:

$$\psi(\mathbf{x}_i, \mathbf{y}) = \sum_{\mathbf{x}_j \in \mathcal{X}_i^+} \sum_{\mathbf{x}_k \in \mathcal{X}_i^-} y_{jk} \left[\frac{\phi(\mathbf{x}_i, \mathbf{x}_j) - \phi(\mathbf{x}_i, \mathbf{x}_k)}{|\mathcal{X}_i^+| \cdot |\mathcal{X}_i^-|} \right]. \quad (4)$$

\mathcal{X}_i^+ and \mathcal{X}_i^- are the sets of relevant and irrelevant neighbours of data point \mathbf{x}_i respectively. Here $|\cdot|$ is the set size. The feature map $\phi(\mathbf{x}_i, \mathbf{x}_j)$ captures the relation between a query \mathbf{x}_i and point \mathbf{x}_j .

We have briefly reviewed how to optimize ranking criteria using structured prediction. Now we review some basic concepts of hashing before introducing our framework.

2.2 Learning-based hashing

Given a set of training data \mathbf{x}_i , ($i = 1, 2, \dots$), the task is to learn a set of hash functions $[h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_\ell(\mathbf{x})]$. Each hash function maps the input into a binary bit $\{0, 1\}$. So with the learned functions, an input \mathbf{x} is mapped into a binary code of length ℓ . We use $\tilde{\mathbf{x}} \in \{0, 1\}^\ell$ to denote the hashed values of \mathbf{x} i.e.,

$$\tilde{\mathbf{x}} = [h_1(\mathbf{x}) \dots, h_\ell(\mathbf{x})]^\top. \quad (5)$$

Suppose that we are given the supervision information as a set of triplets: $\{(\mathbf{x}_i, \mathbf{x}_j, \mathbf{x}_k)\} (i = 1, 2, \dots)$, in which \mathbf{x}_j is an relevant (similar) data point of \mathbf{x}_i (i.e., $\mathbf{x}_j \in \mathcal{X}_i^+$) and \mathbf{x}_k is an irrelevant (dissimilar) neighbor of \mathbf{x}_i (i.e.,

$\mathbf{x}_k \in \mathcal{X}_i^-$). These triplets encode the relative similarity information. After applying the hashing, the distance of two hash codes (of length ℓ) can be calculated using the weighted hamming distance:

$$d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{w}^\top |\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j|, \quad (6)$$

where $\mathbf{w} > 0$ is a non-negative vector, which can be learned. Such weighted hamming distance is used in CGH [12] and multi-dimension spectral hashing [22]. It is expected that after hashing, the distance between relevant data points should be smaller than the distance between irrelevant data points. That is

$$d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_j) \leq d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_k),$$

for $\mathbf{x}_j \in \mathcal{X}_i^+, \mathbf{x}_k \in \mathcal{X}_i^-, \forall i = 1, 2, \dots$. One can then define the margin $\rho = d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_k) - d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_j)$. It is then possible to plug this margin into the large-margin learning framework to optimize for the parameter \mathbf{w} as well as the hash functions, as shown in [12].

3 The proposed StructHash algorithm

For the time being, let us assume that we have already learned all the hashing functions. In other words, *given a data point \mathbf{x} , we assume that we have access to its corresponding hashed values $\tilde{\mathbf{x}}$, as defined in (5)*. Later we will show how this mapping can be explicitly learned using column generation. Now let us focus on how to optimize for the weight \mathbf{w} . When the weighted hamming distance is used, we aim to learn an optimal weight \mathbf{w} defined in (6). Distances are calculated in the learned space and ranked accordingly. A natural choice for the vector-valued mapping function ϕ in Equ. (4) is

$$\phi(\mathbf{x}_i, \mathbf{x}_j) = -|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j|. \quad (7)$$

Note that we have changed the sign, which preserves the ordering in the standard structured SVM. Due to this change of sign, sorting the data by ascending $d_{\text{hm}}(\mathbf{x}_i, \mathbf{x}_j)$ is equivalent to sorting by descending $\mathbf{w}^\top \phi(\mathbf{x}_i, \mathbf{x}_j) = -\mathbf{w}^\top |\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j|$.

The loss function $\Delta(\cdot, \cdot)$ depends on the metric, which we will discuss in detail in the next section. For ease of exposition, let us define

$$\delta\psi_i(\mathbf{y}) = \psi(\mathbf{x}_i, \mathbf{y}_i) - \psi(\mathbf{x}_i, \mathbf{y}). \quad (8)$$

We consider the following problem,

$$\min_{\mathbf{w} \geq 0, \xi \geq 0} \|\mathbf{w}\|_1 + \frac{C}{m} \sum_{i=1}^m \xi_i \quad (9a)$$

$$\text{s.t.: } \forall i = 1, \dots, m \text{ and } \forall \mathbf{y} \in \mathcal{Y} :$$

$$\mathbf{w}^\top \delta\psi_i(\mathbf{y}) \geq \Delta(\mathbf{y}_i, \mathbf{y}) - \xi_i. \quad (9b)$$

Unlike standard structured SVM, here we use the ℓ_1 regularisation (instead of ℓ_2) and enforce that \mathbf{w} to be non-negative. This is aligned with boosting methods [2, 18], and enables us to learn hash functions efficiently.

Algorithm 1 StructHash: Column generation for hash function learning

- 1: **Input:** training examples $(\mathbf{x}_1; \mathbf{y}_1), (\mathbf{x}_2; \mathbf{y}_2), \dots$; parameter C ; the maximum iteration number (bit length ℓ).
 - 2: **Initialise:** working set of hashing functions $\mathcal{W}_H \leftarrow \emptyset$; initialise $\lambda_{(\mathbf{c}, \mathbf{y})} = C/m$ by randomly picking m pairs of (\mathbf{c}, \mathbf{y}) and the rest is set to 0.
 - 3: **Repeat**
 - 4: – Find a new hashing function $h^*(\cdot)$ by solving Equ. (15).
 - 5: – add h^* into the working set of hashing functions: \mathcal{W}_H .
 - 6: – Solve the structured SVM problem (9) or the equivalent (10) using cutting-plane as discussed in Sec. 3.1.
 - 7: **Until** the maximum iteration is reached.
 - 8: **Output:** Learned hash functions and \mathbf{w} .
-

3.1 Learning weights \mathbf{w} via cutting-plane

Here we show how to learn \mathbf{w} . Inspired by [8], we first derive the 1-slack formulation of the original m -slack formulation (9):

$$\min_{\mathbf{w} \geq 0, \xi \geq 0} \|\mathbf{w}\|_1 + C\xi \quad (10a)$$

$$\text{s.t.: } \forall \mathbf{c} \in \{0, 1\}^m \text{ and } \forall \mathbf{y} \in \mathcal{Y}, i = 1, \dots, m:$$

$$\frac{1}{m} \mathbf{w}^\top \left[\sum_{i=1}^m c_i \cdot \delta \psi_i(\mathbf{y}) \right] \geq \frac{1}{m} \sum_{i=1}^m c_i \Delta(\mathbf{y}_i, \mathbf{y}) - \xi. \quad (10b)$$

Here \mathbf{c} enumerates all possible $\mathbf{c} \in \{0, 1\}^m$. As in [8], cutting-plane methods can be used to solve the 1-slack primal problem (10) efficiently. Specifically, we need to solve a maximization for every \mathbf{x}_i in each cutting-plane iteration to find the most violated constraint of (10b), given a solution \mathbf{w} :

$$\mathbf{y}_i^* = \operatorname{argmax}_{\mathbf{y}} \Delta(\mathbf{y}_i, \mathbf{y}) - \mathbf{w}^\top \delta \psi_i(\mathbf{y}). \quad (11)$$

We now know how to efficiently learn \mathbf{w} using cutting-plane methods. However, it remains unclear how to learn hash functions (or features). Thus far, we have taken for granted that the hashed values $\tilde{\mathbf{x}}$ (or h) are given. We would like to learn the hash functions and \mathbf{w} in a single optimization framework. Next we show how this is possible using the column generation technique from boosting.

3.2 Learning hash functions using column generation

Note that the dimension of \mathbf{w} is the same as the dimension of $\tilde{\mathbf{x}}$ (and of $\phi(\cdot, \cdot)$, see Equ. (7)), which is the number of hash bits by the definition (5). If we were able to access all hash functions, it may be possible to select a subset of them and learn the corresponding \mathbf{w} due to the sparsity introduced by the ℓ_1 regularization in (9). Unfortunately, the number of possible hash functions can be infinitely large. In this case it is in general infeasible to solve the optimization problem exactly.

Column generation [2] can be used to approximately solve the problem by adding variables iteratively into the master optimization problems. Column generation was originally invented to solve extremely large-scale linear programming problem, which mainly works on the dual problem. The basic concept of column generation is to add one constraint at a time to the dual problem until an optimal solution is identified. Columns² are generated and added to the problem iteratively to approach the optimality. In the primal problem, column generation solves the problem on a subset of primal variables (\mathbf{w} in our case), which corresponds to a subset of constraints in the dual. This strategy has been widely employed to learn weak learners in boosting [18,2].

To learn hash functions via column generation, we derive the dual problem of the above 1-slack optimization, which is,

$$\max_{\lambda \geq 0} \sum_{\mathbf{c}, \mathbf{y}} \lambda_{(\mathbf{c}, \mathbf{y})} \sum_{i=1}^m c_i \Delta(\mathbf{y}_i, \mathbf{y}) \quad (12a)$$

$$\text{s.t.}: \frac{1}{m} \sum_{\mathbf{c}, \mathbf{y}} \lambda_{(\mathbf{c}, \mathbf{y})} \left[\sum_{i=1}^m c_i \cdot \delta\psi_i(\mathbf{y}) \right] \leq \mathbf{1}, \quad (12b)$$

$$0 \leq \sum_{\mathbf{c}, \mathbf{y}} \lambda_{(\mathbf{c}, \mathbf{y})} \leq C. \quad (12c)$$

We denote by $\lambda_{(\mathbf{c}, \mathbf{y})}$ the 1-slack dual variable associated with one constraint in (10b). Note that (12b) is a set of constraints because $\delta\psi(\cdot)$ is a vector of the same dimension as $\phi(\cdot, \cdot)$ as well as $\tilde{\mathbf{x}}$, which can be infinitely large. One dimension in the vector $\delta\psi(\cdot)$ corresponds to one constraint in (12b). Finding the most violated constraint in the dual form (12) of the 1-slack formulation for generating one hash function is to maximise the l.h.s. of (12b).

The calculation of $\delta\psi(\cdot)$ in (8) can be simplified as follows. Because of the subtraction of $\psi(\cdot)$ (defined in (4)), only those incorrect ranking pairs will appear in the calculation. Recall that the true ranking is \mathbf{y}_i for \mathbf{x}_i . We define $\mathcal{S}_i(\mathbf{y})$ as a set of incorrectly ranked pairs: $(j, k) \in \mathcal{S}_i(\mathbf{y})$, in which the incorrectly ranked pair (j, k) means that the true ranking is $\mathbf{x}_j \prec_{\mathbf{y}_i} \mathbf{x}_k$ but $\mathbf{x}_j \succ_{\mathbf{y}} \mathbf{x}_k$. So we have

$$\begin{aligned} \delta\psi_i(\mathbf{y}) &= \frac{2}{|\mathbf{x}_i^+| |\mathbf{x}_i^-|} \sum_{(j, k) \in \mathcal{S}_i(\mathbf{y})} [\phi(\mathbf{x}_i, \mathbf{x}_j) - \phi(\mathbf{x}_i, \mathbf{x}_k)] \\ &= \frac{2}{|\mathbf{x}_i^+| |\mathbf{x}_i^-|} \sum_{(j, k) \in \mathcal{S}_i(\mathbf{y})} (|\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_k| - |\tilde{\mathbf{x}}_i - \tilde{\mathbf{x}}_j|). \end{aligned} \quad (13)$$

² A column is a variable in the primal and a corresponding constraint in the dual.

With the above equations and the definition of $\tilde{\mathbf{x}}$ in (5), the most violated constraint in (12b) can be found by solving the following problem:

$$h^*(\cdot) = \underset{h(\cdot)}{\operatorname{argmax}} \sum_{\mathbf{c}, \mathbf{y}} \lambda_{(\mathbf{c}, \mathbf{y})} \sum_i \frac{2c_i}{|\mathcal{X}_i^+| |\mathcal{X}_i^-|} \cdot \sum_{(j,k) \in \mathcal{S}_i(\mathbf{y})} (|h(\mathbf{x}_i) - h(\mathbf{x}_k)| - |h(\mathbf{x}_i) - h(\mathbf{x}_j)|). \quad (14)$$

By exchanging the order of summations, the above optimization can be further written in a compact form:

$$h^*(\cdot) = \underset{h(\cdot)}{\operatorname{argmax}} \sum_{i, \mathbf{y}} \sum_{(j,k) \in \mathcal{S}_i(\mathbf{y})} \mu_{(i, \mathbf{y})} (|h(\mathbf{x}_i) - h(\mathbf{x}_k)| - |h(\mathbf{x}_i) - h(\mathbf{x}_j)|), \quad (15)$$

$$\text{where, } \mu_{(i, \mathbf{y})} = \frac{2}{|\mathcal{X}_i^+| |\mathcal{X}_i^-|} \sum_{\mathbf{c}} \lambda_{(\mathbf{c}, \mathbf{y})} c_i. \quad (16)$$

The objective in the above optimization is a summation of weighted triplet (i, j, k) ranking scores, in which $\mu_{(i, \mathbf{y})}$ is the triplet weighting value. Solving the above optimization provides the best hash function for the current solution \mathbf{w} . Once a hash function is generated, we learn \mathbf{w} using cutting-plane in Sec. 3.1. The column generation procedure for hash function learning is summarised in Algorithm 1.

The form of hash function $h(\cdot)$ can be any function that have binary output value. For a decision stump as the hash function, usually we can exhaustively enumerate all possibility and find the globally best one. However globally solving (15) is generally difficult. In our experiments, we use the linear perceptron hash function with the output in $\{0, 1\}$:

$$h(\mathbf{x}) = 0.5(\operatorname{sign}(\mathbf{v}^\top \mathbf{x} + b) + 1). \quad (17)$$

The non-smooth function $\operatorname{sign}(\cdot)$ here brings the difficulty for optimization. Similar to [12], we replace the $\operatorname{sign}(\cdot)$ function by a smooth sigmoid function, and then locally solve the above optimization (15) (e.g., LBFGS [26]) for learning the parameters of a hash function. We can apply a few heuristics to initialize for solving (15). For example, similar to LSH, we can generate a set of random projection planes and then choose the best one that maximizes the objective in (15) as the initialization. We can also train a decision stump by searching a best dimension and threshold to maximize the objective on the quantized data. Alternatively, one can employ the spectral relaxation method [16] which drops the $\operatorname{sign}(\cdot)$ function and solves a generalized eigenvalue problem to obtain an initial point. In our experiments, we use the spectral relaxation method for initialization.

Next, we discuss some widely-used information retrieval evaluation criteria, and show how they can be seamlessly incorporated into StructHash.

4 Ranking measures

Here we discuss a few ranking measures for loss functions, including AUC, NDCG, Precision-at-K, and mAP. Following [17], we define the loss function over two rankings $\Delta \in [0, 1]$ as:

$$\Delta(\mathbf{y}, \mathbf{y}') = 1 - \text{score}(\mathbf{y}, \mathbf{y}'). \quad (18)$$

Here \mathbf{y}' is the ground truth ranking and \mathbf{y} is the prediction. We define $\mathcal{X}_{\mathbf{y}'}^+$ and $\mathcal{X}_{\mathbf{y}'}^-$ as the indexes of relevant and irrelevant neighbours respectively in the ground truth ranking \mathbf{y}' .

AUC. The area under the ROC curve is to evaluate the performance of correct ordering of data pairs, which can be computed by counting the proportion of correctly ordered data pairs:

$$\text{score}_{\text{AUC}}(\mathbf{y}, \mathbf{y}') = \frac{1}{|\mathcal{X}_{\mathbf{y}'}^+| |\mathcal{X}_{\mathbf{y}'}^-|} \sum_{i \in \mathcal{X}_{\mathbf{y}'}^+} \sum_{j \in \mathcal{X}_{\mathbf{y}'}^-} \delta(i \prec_{\mathbf{y}} j). \quad (19)$$

$\delta(\cdot) \in \{0, 1\}$ is the indicator function. For using this AUC loss, the maximization inference in (11) can be solved efficiently by sorting the distances of data pairs, as described in [7]. Note that the loss of a wrongly ordered pair is not related to their positions in the ranking list, thus AUC is a position insensitive measure.

Precision-at-K. Precision-at-K is to evaluate the quality of top-K retrieved examples in a ranking. It is computed by counting the number of relevant data points within top-K positions and divided by K:

$$\text{score}_{\text{P@K}}(\mathbf{y}, \mathbf{y}') = \frac{1}{K} \sum_{i=1}^K \delta(\mathbf{y}(i) \in \mathcal{X}^+). \quad (20)$$

Here $\mathbf{y}(i)$ is the example index on the i -th position of a ranking \mathbf{y} ; $\delta(\cdot)$ is an indicator. An algorithm for solving the inference in (11) is proposed in [7].

NDCG. Normalized Discounted Cumulative Gain [6] is to measure the ranking quality of the first K returned neighbours. A similar measure is Precision-at-K which is the proportion of top-K relevant neighbours. NDCG is a position-sensitive measure which considers the positions of the top-K relevant neighbours. Compared to the position-insensitive measure: AUC, NDCG assigns different importances on the ranking positions, which is a more favorable measure for a general notion of a “good” ranking in real-world applications. In NDCG, each position of the ranking is assigned a score in a decreasing way. NDCG can be computed by accumulating the scores of top-K relevant neighbours:

$$\text{score}_{\text{NDCG}}(\mathbf{y}, \mathbf{y}') = \frac{1}{\sum_{i=1}^K S(i)} \sum_{i=1}^K S(i) \delta(\mathbf{y}(i) \in \mathcal{X}_{\mathbf{y}'}^+). \quad (21)$$

Here $\mathbf{y}(i)$ is the example index on the i -th position of a ranking \mathbf{y} . $S(i)$ is the score assigned to the i -th position of a ranking. $S(1) = 1$, $S(i) = 0$ for $i > k$ and

Table 1. Results using NDCG measure (64 bits). We compare our StructHash using AUC (StructH-A) and NDCG (StructH-N) loss functions with other supervised and un-supervised methods. Our method using NDCG loss performs the best in most cases.

Dataset	StructH-N	StructH-A	CGH	SPLH	STHs	BREs	ITQ	SPHER	MDSH	AGH	LSH
NDCG ($K = 100$)											
STL10	0.435	0.374	0.375	0.404	0.214	0.289	0.337	0.318	0.313	0.310	0.228
USPS	0.905	0.893	0.900	0.816	0.688	0.777	0.804	0.762	0.735	0.741	0.668
MNIST	0.851	0.798	0.867	0.804	0.594	0.805	0.856	0.806	0.100	0.793	0.561
CIFAR	0.335	0.259	0.258	0.357	0.178	0.273	0.314	0.297	0.283	0.286	0.168
ISOLET	0.881	0.839	0.866	0.629	0.766	0.483	0.623	0.518	0.538	0.536	0.404

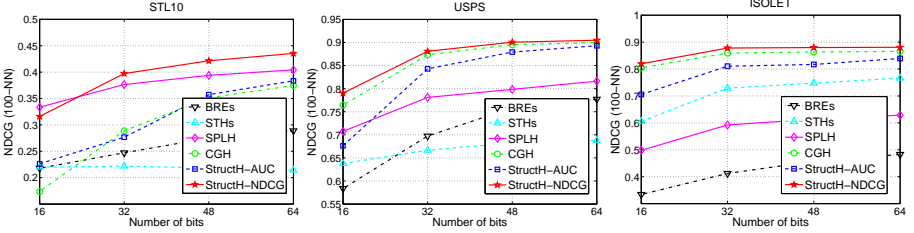


Fig. 1. NDCG results on 3 datasets. Our StructHash performs the best.

$S(i) = 1/\log_2(i)$ for other cases. A dynamic programming algorithm is proposed in [1] for solving the maximization inference in (11).

mAP. Mean average precision (mAP) is the averaged precision-at-K scores over all positions of relevant data points in a ranking, which is computed as:

$$\text{score}_{\text{mAP}}(\mathbf{y}, \mathbf{y}') = \frac{1}{|\mathcal{X}_{\mathbf{y}'}^+|} \sum_{i=1}^{|\mathcal{X}_{\mathbf{y}'}^+| + |\mathcal{X}_{\mathbf{y}'}^-|} \delta(\mathbf{y}(i) \in \mathcal{X}_{\mathbf{y}'}^+) \text{score}_{\text{P@K}(K=i)}(\mathbf{y}, \mathbf{y}'). \quad (22)$$

For using this mAP loss, an efficient algorithm for solving the inference in (11) is proposed in [24].

5 Experiments

Our method is in the category of supervised method for learning compact binary codes. Thus we mainly compare with 4 supervised methods: column generation hashing (CGH) [12], supervised binary reconstructive embeddings (BREs) [10], supervised self-taught hashing (STHs) [25], semi-supervised sequential projection learning hashing (SPLH) [21].

For comparison, we also run some unsupervised methods: locality-sensitive hashing (LSH) [3], anchor graph hashing (AGH) [16], spherical hashing (SPHER) [5], multi-dimension spectral hashing (MDSH) [22], and iterative quantization (ITQ) [4]. We carefully follow the original authors' instruction for parameter setting. For SPLH, the regularization parameter is picked from 0.01 to 1. We use the hierarchical variant of AGH. The bandwidth parameters of Gaussian affinity



Fig. 2. Some ranking examples of our method. The first column shows query images, and the rest are retrieved images. False predictions are marked by red boxes.

in MDSH is set as $\sigma = t\bar{d}$. Here \bar{d} is the average Euclidean distance of top 100 nearest neighbours and t is picked from 0.01 to 50. For supervised training of our StructHash and CGH, we use 50 relevant and 50 irrelevant examples to construct similarity information for each data point.

We use 9 datasets for evaluation, including one UCI dataset: ISOLET, 4 image datasets: CIFAR10³, STL10⁴, MNIST, USPS, and another 4 large image datasets: Tiny-580K[4], Flickr-1M⁵, SIFT-1M [21] and GIST-1M⁶. CIFAR10 is a subset of the 80-million tiny images and STL10 is a subset of Image-Net. Tiny-580K consists of 580,000 tiny images. Flickr-1M dataset consists of 1 million thumbnail images. SIFT-1M and GIST-1M datasets contain 1 million SIFT and GIST features respectively.

We follow a common setting in many supervised methods [10,15,12] for hashing evaluation. For multi-class datasets, we use class labels to define the relevant and irrelevant semantic neighbours by label agreement. For large datasets: Flickr-1M, SIFT-1M, GIST-1M and Tiny-580K, the semantic ground truth is

³ <http://www.cs.toronto.edu/~kriz/cifar.html>

⁴ <http://www.stanford.edu/~acoates/stl10/>

⁵ <http://press.liacs.nl/mirflickr/>

⁶ <http://corpus-texmex.irisa.fr/>

Table 2. Results using ranking measures of Precision-at-K, Mean Average Precision and Precision-Recall (64 bits). We compare our method using AUC (StructH-A) and NDCG (StructH-N) loss functions with other supervised and un-supervised methods. Our method using NDCG loss performs the best on these measures

Dataset	StructH-N	StructH-A	CGH	SPLH	STHs	BREs	ITQ	SPHER	MDSH	AGH	LSH
Precision-at-K ($K = 100$)											
STL10	0.431	0.376	0.376	0.396	0.208	0.279	0.325	0.303	0.298	0.301	0.222
USPS	0.903	0.894	0.898	0.805	0.667	0.755	0.780	0.730	0.698	0.711	0.637
MNIST	0.849	0.807	0.862	0.797	0.579	0.790	0.842	0.788	0.100	0.780	0.540
CIFAR	0.336	0.259	0.261	0.354	0.174	0.264	0.301	0.286	0.270	0.281	0.164
ISOLET	0.875	0.844	0.859	0.604	0.755	0.448	0.589	0.477	0.493	0.493	0.370
Mean Average Precision (mAP)											
STL10	0.331	0.326	0.322	0.299	0.155	0.211	0.233	0.193	0.178	0.162	0.162
USPS	0.868	0.851	0.848	0.689	0.456	0.582	0.566	0.451	0.405	0.333	0.418
MNIST	0.802	0.790	0.789	0.684	0.397	0.558	0.585	0.510	0.119	0.505	0.343
CIFAR	0.294	0.300	0.298	0.289	0.147	0.204	0.215	0.204	0.181	0.201	0.149
ISOLET	0.836	0.796	0.815	0.518	0.653	0.340	0.484	0.357	0.348	0.298	0.267
Precision-Recall											
STL10	0.267	0.248	0.248	0.246	0.130	0.181	0.200	0.174	0.164	0.145	0.138
USPS	0.776	0.760	0.760	0.609	0.401	0.520	0.508	0.424	0.379	0.326	0.375
MNIST	0.591	0.574	0.582	0.445	0.165	0.313	0.323	0.246	0.018	0.197	0.143
CIFAR	0.105	0.093	0.091	0.110	0.042	0.066	0.074	0.069	0.064	0.061	0.042
ISOLET	0.759	0.709	0.737	0.445	0.563	0.301	0.429	0.321	0.320	0.275	0.238

defined according to the ℓ_2 distance [21]. Specifically, a data point is labeled as a relevant data point of the query if it lies in the top 2 percentile points in the whole dataset. We generated GIST features for all image datasets except MNIST and USPS. we randomly select 2000 examples for testing queries, and the rest is used as database. We sample 2000 examples from the database as training data for learning models. For large datasets, we use 5000 examples for training. To evaluate the performance of compact bits, the maximum bit length is set to 64, as similar to the evaluation settings in other supervised hashing methods [10,12].

We report the result of the NDCG measure in Table 1. We compare our StructHash using AUC and NDCG loss functions with other supervised and un-supervised methods. Our method using NDCG loss function performs the best in most cases. We also report the result of other common measures in Table 2, including the result of Precision-at-K, Mean Average Precision (mAP) and Precision-Recall. Precision-at-K is the proportion of true relevant data points in the returned top-K results. The Precision-Recall curve measures the overall performance in all positions of the prediction ranking, which is computed by varying the number of nearest neighbours. It shows that our method generally performs better than other methods on these evaluation measures. As described before, compared to the AUC measure which is position insensitive, the NDCG measure assigns different importance on ranking positions, which is closely related to many other position sensitive ranking measures (e.g., mAP). As expected, the result shows that on the Precision-at-K, mAP and Precision-recall measures, optimizing the position sensitive NDCG loss performs better than the AUC loss. The triplet loss based method CGH actually helps to reduce the AUC loss. This may explain the reason that our method, which aims to reduce the NDCG loss,

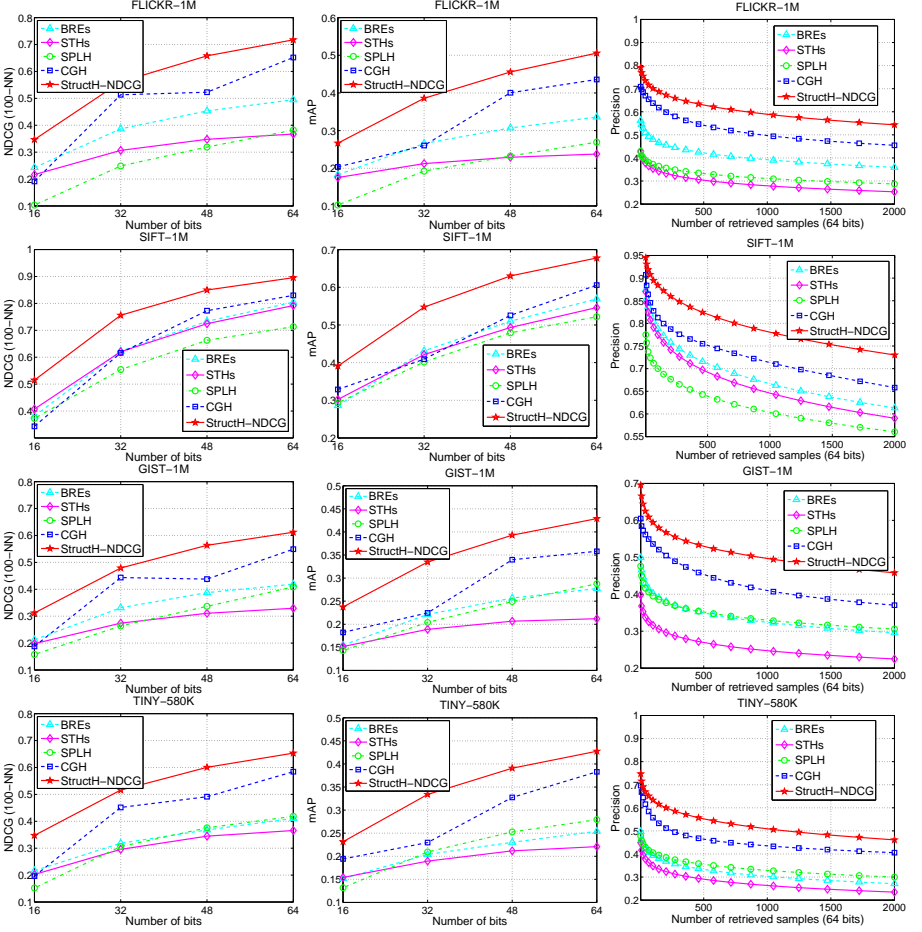


Fig. 3. Results on 4 large datasets: Flickr-1M (1 million Flickr images), Sift-1M (1 million SIFT features), Gist-1M (1 million GIST features) and Tiny580K (580,000 Tiny image dataset). We compare with several supervised methods. The results of 3 measures (NDCG, mAP and precision of top-K neighbours) are shown here. Our StructHash outperforms others in most cases.

is able to outperform CGH in these measures. We also plot the NDCG results on several datasets in Fig. 1 by varying the number of bits. Some retrieval examples are shown in Fig. 2.

We further evaluate our method on 4 large-scale datasets (Flickr-1M, SIFT-1M, GIST-1M and Tiny-580K). The results of NDCG, mAP and the precision of top-K neighbours are shown in Fig. 3. The NDCG and mAP results are shown by varying the number of bits. The precision of top-K neighbours is shown by varying the number of retrieved examples. In most cases, our method outperforms

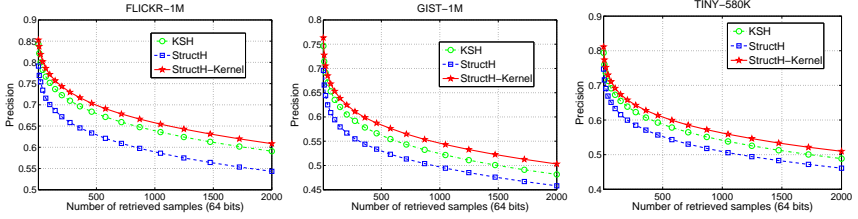


Fig. 4. Comparison on large datasets of our kernel StructHash (StructHash-Kernel) with our non-kernel StructHash and the relevant method KSH [15]. Our kernel version is able to achieve better results.

other competitors. Our method with NDCG loss function succeeds to achieve good performance both on NDCG and other measures.

Applying the kernel technique in KLSH [11] and KSH [15] further improves the performance of our method. As describe in [15], we perform a pre-processing step to generate the kernel mapping features: we randomly select a number of support vectors (300) then compute the kernel response on data points as input features. Note that here we simply follow KSH for the kernel parameter setting. We evaluate this kernel version of our method in Fig. 4 and compare to KSH. Our kernel version is able to achieve better results.

6 Conclusion and future work

We have developed a hashing framework that allows us to directly optimize multivariate performance measures. The fact that the proposed method outperforms comparable hashing approaches is to be expected, as it more directly optimizes the required loss function. It is anticipated that the success of the approach may lead to a range of new hashing-based applications with task-specified targets. Extension to make use of the power of more sophisticated hash functions such as kernel functions or decision trees is of future work.

References

1. Chakrabarti, S., Khanna, R., Sawant, U., Bhattacharyya, C.: Structured learning for non-smooth ranking losses. In: Proc. ACM Knowledge Discovery & Data Mining (2008)
2. Demiriz, A., Bennett, K.P., Shawe-Taylor, J.: Linear programming boosting via column generation. Mach. Learn. (2002)
3. Gionis, A., Indyk, P., Motwani, R.: Similarity search in high dimensions via hashing. In: Proc. Int. Conf. Very Large Data Bases (1999)
4. Gong, Y., Lazebnik, S., Gordo, A., Perronnin, F.: Iterative quantization: a procrustean approach to learning binary codes for large-scale image retrieval. IEEE Trans. Patt. Anal. & Mach. Intelli. (2012)
5. Heo, J., Lee, Y., He, J., Chang, S., Yoon, S.: Spherical hashing. In: Proc. Int. Conf. Comp. Vis. & Patt. Recogn. (2012)

6. Järvelin, K., Kekäläinen, J.: IR evaluation methods for retrieving highly relevant documents. In: Proc. ACM Conf. SIGIR (2000)
7. Joachims, T.: A support vector method for multivariate performance measures. In: Proc. Int. Conf. Mach. Learn. (2005)
8. Joachims, T.: Training linear SVMs in linear time. In: Proc. ACM Knowledge Discovery & Data Mining (2006)
9. Kelley, J. E., J.: The cutting-plane method for solving convex programs. J. Society for Industrial & Applied Math. (1960)
10. Kulis, B., Darrell, T.: Learning to hash with binary reconstructive embeddings. Proc. Adv. Neural Info. Process. Syst. (2009)
11. Kulis, B., Grauman, K.: Kernelized locality-sensitive hashing. IEEE Trans. Patt. Anal. & Mach. Intelli. (2012)
12. Li, X., Lin, G., Shen, C., van den Hengel, A., Dick, A.: Learning hash functions using column generation. In: Proc. Int. Conf. Mach. Learn. (2013)
13. Lin, G., Shen, C., Shi, Q., van den Hengel, A., Suter, D.: Fast supervised hashing with decision trees for high-dimensional data. In: Proc. Int. Conf. Comp. Vis. & Patt. Recogn. Columbus, Ohio, USA (2014), <https://bitbucket.org/chhshen/fasthash/src>
14. Lin, G., Shen, C., Suter, D., van den Hengel, A.: A general two-step approach to learning-based hashing. In: Proc. Int. Conf. Comp. Vis. Sydney, Australia (2013)
15. Liu, W., Wang, J., Ji, R., Jiang, Y., Chang, S.: Supervised hashing with kernels. In: Proc. Int. Conf. Comp. Vis. & Patt. Recogn. (2012)
16. Liu, W., Wang, J., Kumar, S., Chang, S.F.: Hashing with graphs. In: Proc. Int. Conf. Mach. Learn. (2011)
17. McFee, B., Lanckriet, G.: Metric learning to rank. In: Proc. Int. Conf. Mach. Learn. (2010)
18. Shen, C., Li, H.: On the dual formulation of boosting algorithms. IEEE Trans. Patt. Anal. & Mach. Intelli. (2010)
19. Shen, F., Shen, C., Shi, Q., van den Hengel, A., Tang, Z.: Inductive hashing on manifolds. In: Proc. Int. Conf. Comp. Vis. & Patt. Recogn. Oregon, USA (2013)
20. Tsochantaridis, I., Hofmann, T., Joachims, T., Altun, Y.: Support vector machine learning for interdependent and structured output spaces. In: Proc. Int. Conf. Mach. Learn. (2004)
21. Wang, J., Kumar, S., Chang, S.: Semi-supervised hashing for large scale search. IEEE Trans. Patt. Anal. & Mach. Intelli. (2012)
22. Weiss, Y., Fergus, R., Torralba, A.: Multidimensional spectral hashing. In: Proc. Eur. Conf. Comp. Vis. (2012)
23. Weiss, Y., Torralba, A., Fergus, R.: Spectral hashing. In: Proc. Adv. Neural Info. Process. Syst. (2008)
24. Yue, Y., Finley, T., Radlinski, F., Joachims, T.: A support vector method for optimizing average precision. In: Proc. ACM Conf. SIGIR (2007)
25. Zhang, D., Wang, J., Cai, D., Lu, J.: Extensions to self-taught hashing: kernelisation and supervision. In: Proc. ACM Conf. SIGIR Workshop (2010)
26. Zhu, C., Byrd, R.H., Lu, P., Nocedal, J.: Algorithm 778: L-BFGS-B: Fortran sub-routines for large-scale bound-constrained optimization. ACM T. Math. Softw. (1997)